

# REPORT DOCUMENTATION PAGE

AFRL-SR-AR-TR-04-

Public reporting burden for this collection of information is estimated to average 1 hour per response, including the time for reviewing existing data needed, and completing and reviewing this collection of information. Send comments regarding this burden estimate suggestions for reducing this burden to Washington Headquarters Services, Directorate for Information Operations and Reports, 1215 Jefferson Davis Highway, Suite 1204, Arlington, VA 22202-4302, and to the Office of Management and Budget, Paperwork Reduction Project (0704-0188), Washington, DC 20503

0162

1. AGENCY USE ONLY (Leave blank)		2. REPORT DATE: 3/04		3. REPORT TYPE AND DATES COVERED Final technical report	
4. TITLE AND SUBTITLE  Computational Auditory Scene Analysis Based Perceptual and Neural Principles				5. FUNDING NUMBERS  F49620-01-1-0027	
6. AUTHOR(S)  DeLiang Wang: Principal Investigator					
7. PERFORMING ORGANIZATION NAME(S) AND ADDRESS(ES)  The Ohio State University Research Foundation 1960 Kenny Road Columbus, OH 43210-1063				8. PERFORMING ORGANIZATION REPORT NUMBER 740381/868490	
9. SPONSORING / MONITORING AGENCY NAME(S) AND ADDRESS(ES)  AFOSR Attn: Dr. Willard Larkin AFOSR/NL, Room 713 4015 Wilson Blvd. Arlington, VA 22203-1945				10. SPONSORING / MONITORING AGENCY REPORT NUMBER	
11. SUPPLEMENTARY NOTES					
12a. DISTRIBUTION / AVAILABILITY STATEMENT Distribution unrestricted				20040319 117	
13. ABSTRACT (Maximum 200 Words) A remarkable feat of the auditory system is its ability to disentangle the acoustic mixture and group the acoustic energy from the same event. This fundamental process of auditory perception is called auditory scene analysis. Of particular importance in auditory scene analysis is the separation of speech from interfering sounds, or speech segregation. Consistent with specified objectives, this project made major advances along the following three directions. First, the problem of multipitch tracking was investigated in the context of multiple sound sources, and a robust algorithm for multipitch tracking of noisy speech was developed. The second advance is in monaural separation of voiced speech, where a new system was proposed that employs different strategies in the low- and the high-frequency range. A key element of the system is amplitude modulation analysis in the high-frequency range. Third, the problem of location-based separation was studied in the joint feature space of interaural time difference and interaural intensity difference, and a novel classification approach was introduced to optimally determine whether a target sound dominates in local time-frequency units. All of the three models were comprehensively evaluated and shown to be substantially superior to existing approaches.					
14. SUBJECT TERMS Auditory scene analysis, speech segregation, multipitch tracking, amplitude modulation, binaural segregation, computational audition				15. NUMBER OF PAGES 13	
				16. PRICE CODE	
17. SECURITY CLASSIFICATION OF REPORT Unclassified	18. SECURITY CLASSIFICATION OF THIS PAGE Unclassified	19. SECURITY CLASSIFICATION OF ABSTRACT Unclassified	20. LIMITATION OF ABSTRACT		

NSN 7540-01-280-5500

Standard Form 298 (Rev. 2-89)  
Prescribed by ANSI Std. Z39-18  
298-102

## DISTRIBUTION STATEMENT A

Approved for Public Release  
Distribution Unlimited

# Final Performance Report

DeLiang Wang

*The Ohio State University*

This PI was awarded an AFOSR grant for the project "Computational auditory scene analysis based on perceptual and neural principles" (Grant No. F49620-01-1-0027). The project was funded for the period of 12/1/00 through 11/30/03 with the total amount \$486K. This report summarizes the progress made throughout the project period.

## 1. RESEARCH PROGRESS

Auditory scene analysis is the perceptual process in which sounds from different sources are separated into individual representations. Consistent with specified objectives, the project has made major advances along the following three directions. First, we have worked on the problem of multipitch tracking in the context of multiple sound sources, and have developed a robust algorithm for multipitch tracking of noisy speech. Second, we have investigated the problem of monaural separation of voiced speech, and the resulting model performs significantly better than previous systems. Third, we have studied the problem of binaural sound separation using joint ITD (interaural time difference) and IID (interaural intensity difference) cues, and have introduced a novel classification approach to optimally determine whether a target sound dominates in a local time-frequency unit. The major findings along these three directions are described in the following subsections.

### 1.1 Multipitch Tracking in Noisy Environments

A reliable algorithm for tracking multipitch contours is needed for not only sound source separation, but also prosody analysis and speaker recognition. However, due to the difficulty of dealing with noise intrusions and mutual interference among multiple harmonic structures, the design of such an algorithm has proven to be very challenging. Most existing pitch determination algorithms ignore the multipitch nature of an acoustic mixture, and are limited to clean speech or a single pitch track with modest background noise. To be useful for computational auditory scene analysis (CASA), a pitch determination algorithm must perform in a variety of acoustic environments.

We have developed a robust algorithm for multipitch tracking of noisy speech. An input mixture is first processed by a bank of gammatone filters, which closely model the cochlear filtering process. Our pitch-tracking algorithm builds on the correlogram representation, which computes autocorrelation functions of gammatone filter responses. Correlograms provide a joint time-frequency representation. Our algorithm combines an improved method for channel and peak selection and a new method for extracting and integrating periodicity across different filter channels. The algorithm maintains multiple

hypotheses with different probabilities, and the modeling process incorporates the statistics extracted from natural sound sources. Finally, a hidden Markov model (HMM) is employed to detect continuous pitch tracks.

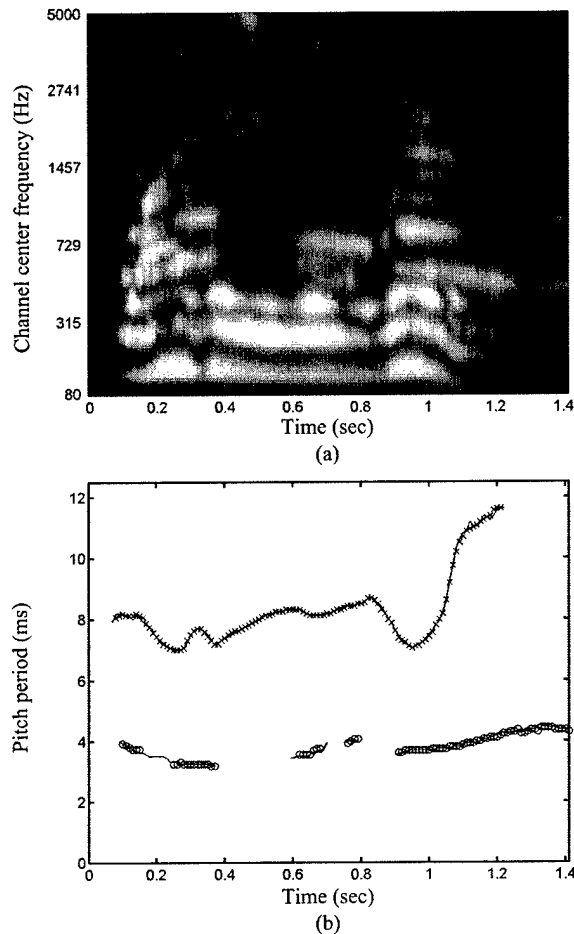
The resulting algorithm has been comprehensively evaluated and compared with existing algorithms. The evaluation results show that the algorithm can reliably extract single and double pitch tracks in a noisy environment and outperforms other algorithms by a large margin. Figure 1 shows an example. Figure 1a shows the time-frequency energy plot for a mixture of two simultaneous utterances: A male speaker and a female speaker with an energy ratio of 9 dB. Figure 1b shows the multipitch tracking result generated by the algorithm.

A preliminary version of the work was presented in the 2002 *International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, and the complete version has been recently published by *IEEE Transactions on Speech and Audio Processing*. Detailed references are given in Section 2.3 (Wu, Wang, and Brown, 2002; 2003). Also, a powerpoint presentation file - titled "Multipitch tracking for noisy speech" - is posted on the PI's webpage at <http://www.cis.ohio-state.edu/~dwang> (follow the Presentations link), and the file contains sound demos. In addition, the source program that implements the algorithm is posted at the PI's laboratory webpage at <http://www.cis.ohio-state.edu/pnl> (follow the Software link).

## 1.2 Monaural Segregation of Voiced Speech

For voiced speech, harmonicity is an essential cue for segregation. Former CASA systems are able to segregate most low-frequency voiced speech signals, but are incapable of segregating high-frequency speech signals. It is well-known that the auditory system can resolve the first few harmonics, while higher harmonics are unresolved. Psychoacoustic evidence suggests that the auditory system uses different mechanisms to deal with resolved and unresolved harmonics. Existing CASA systems use the same method to deal with both low-frequency and high-frequency signals, which is a main reason that they have difficulty dealing with high-frequency responses.

We have developed a system that employs different methods in the low- and the high-frequency range. In the low-frequency range, our system generates segments based on temporal continuity and common periodicity between responses of adjacent channels. These segments are grouped by comparing their periodicities with estimated pitch periods of the target speech. In the high-frequency range, wide bandwidths of auditory filters make the filters respond to multiple unresolved harmonics of voiced speech. These responses are amplitude modulated and their envelopes fluctuate at the frequency that corresponds to the fundamental frequency ( $F_0$ ). Hence, our system generates segments in the high-frequency range based on temporal continuity and common amplitude modulation (AM) among adjacent filter responses. These segments are grouped by comparing AM repetition rates with estimated  $F_0$  of target speech. To derive AM repetition rates we have employed a sinusoidal modeling technique; specifically, we use a single sinusoid to model AM repetition within a certain pitch range, and the derivation of AM repetition rates can then be formulated as an optimization problem. With appropriately chosen initial values, the optimization problem can be solved efficiently using an iterative gradient descent technique. In addition, we have developed



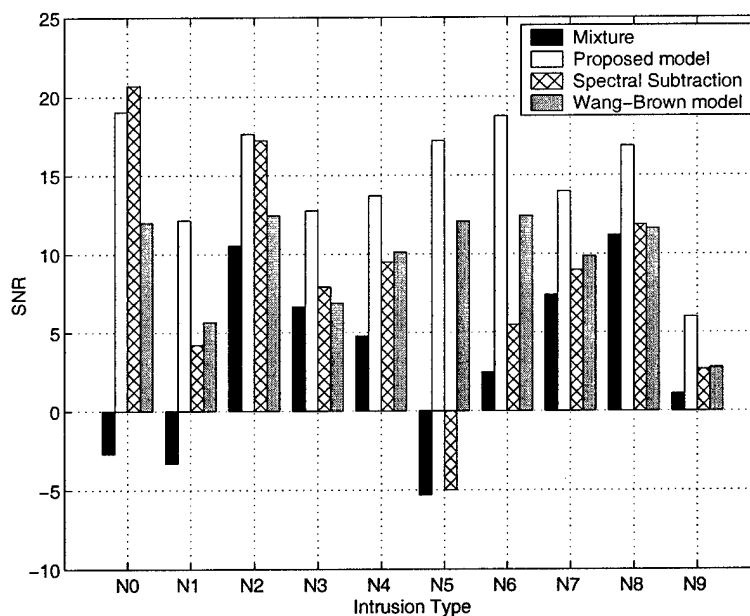
**Figure 1.** (a) Time-frequency energy plot for a mixture of two simultaneous utterances of a male and a female speaker. The male utterance is “Why are you all weary” and the female utterance is “Don’t ask me to carry an oily rag like that.” The brightness in a time-frequency unit (pixel) indicates the energy of the corresponding gammatone filter output in the corresponding time frame. (b) Result of tracking the mixture. The solid lines indicate the true (correct) pitch tracks. The ‘x’ and ‘o’ tracks represent the pitch tracks estimated by the proposed algorithm.

a psychoacoustically motivated method for tracking a single target pitch contour (unlike multipitch tracking described in Section 1.1).

The above analysis of amplitude modulation in the high-frequency range, along with a pitch contour from target tracking, yields a monaural segregation system for voiced speech. The system has been systematically evaluated, and it produces substantially better performance than previous models, especially in the high-frequency range. Figure 2 shows that, using the a mixture database of voiced speech and interference, the resulting

system produces significant improvement in signal-to-noise ratio (SNR) over that of the Wang and Brown model published in 1999, which had the representative performance of previous CASA systems, and the standard spectral subtraction method in speech enhancement. The figure shows the SNR for each intrusion averaged across 10 target utterances, together with the SNR of the original mixtures and the results from the Wang-Brown model and spectral subtraction. All three systems show SNR improvements over original mixtures. Compared to the Wang-Brown model, the new system yields at least 3 dB SNR improvement for every intrusion type. The average improvement for the entire corpus is about 5.2 dB. The Wang-Brown model in turn performs 1.2 dB better on average than spectral subtraction. As expected, spectral subtraction produces uneven results for the intrusions; for example, its performance is the best among all the methods for the pure tone intrusion (N0).

Preliminary reports describing this work were presented in 2002 *ICASSP* and the 2002 *Neural Information Processing Systems (NIPS) Conference*. An extensive version was recently accepted by *IEEE Transactions on Neural Networks* (a hardcopy is attached). Detailed references are given in Section 2.3 (Hu and Wang, 2002a; 2002b; 2004). Also, a powerpoint presentation file - titled "Monaural speech segregation" - is posted on the PI's webpage at <http://www.cis.ohio-state.edu/~dwang> (follow the Presentations link), and the file contains sound demos. In addition, the source program that implements the algorithm is posted at the PI's laboratory webpage at <http://www.cis.ohio-state.edu/pnl> (follow the Software link).



**Figure 2.** SNR results for segregated speech and original mixtures. White bars show the results from the proposed model, gray bars from the Wang-Brown system, cross bars from the spectral subtraction method, and black bars from original mixtures. The intrusion types are N0: 1 kHz tone; N1: random noise; N2: noise burst; N3: 'cocktail party'; N4: rock music; N5: siren; N6: telephone; N7: female speech; N8: male speech; and N9: female speech.

### 1.3 Location-based Speech Segregation

At a cocktail party, people can selectively attend to a single voice and filter out other acoustical interferences, and location plays an important role in cocktail party processing. How to simulate this perceptual ability, known as the cocktail-party problem, remains a great challenge.

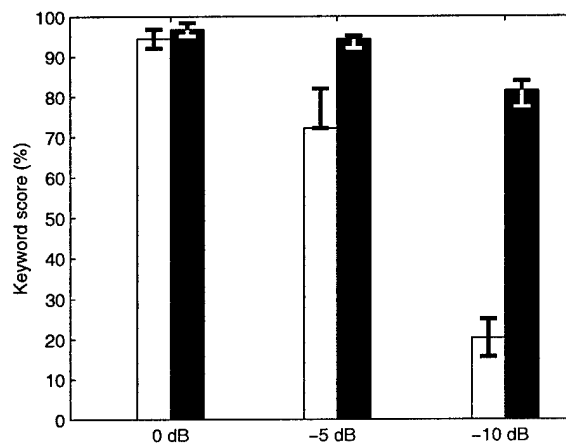
We have developed a novel location-based approach to speech segregation. Our model starts with the binaural cues of ITD and IID extracted from the responses of a KEMAR dummy head that realistically simulates the filtering process of the head, torso and external ear. Motivated by the auditory masking phenomenon, we have introduced the notion of an “ideal” time-frequency binary mask. The ideal binary mask selects the target if it is stronger than the interference in a local time-frequency unit. We observe that, within a narrow frequency band, modifications to the relative energy of the target source with respect to the interfering energy trigger systematic changes of the binaural cues. For a given spatial configuration, this interaction produces characteristic clustering in the binaural feature space. Consequently, we employ a classification method to determine decision regions in the joint ITD-IID feature space that correspond to target estimates.

Our system has been systematically evaluated for two-source configurations where the target position moves from the median plane to the side of the head. Excellent results are obtained for target in the median plane for an azimuth separation as small as  $5^\circ$ . Performance degrades when the target source is moved to the side of the head, where a  $10^\circ$  separation is needed for good performance. This is a direct consequence of poorer resolution of the binaural cues to the side. This performance profile is in agreement with empirical data from human observers. When comparing the SNR with that of the initial mixture, there is an average SNR gain of 14 dB for target sources in the median plane. This reduces to 11 dB when the target source is at  $70^\circ$ . For 3 sources, average SNR gain is 11.3 dB in favorable configurations. We have also implemented and compared with the Bodden model published in 1993, which estimates a Wiener filter for speech segregation. Our system produces 3.5 dB improvement in the most favorable conditions for the Bodden model, and in other configuration conditions the improvement is significantly greater.

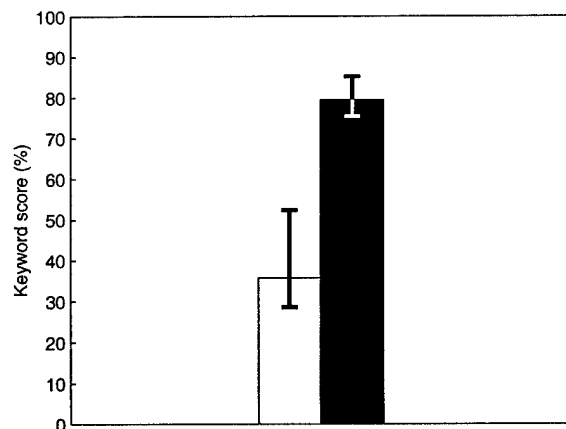
In addition to SNR evaluation, we have also performed a systematic evaluation in terms of automatic speech recognition (ASR) performance, and the results confirm that our system produces masks very close to ideal binary ones and yields large ASR improvements compared to direct recognition of mixtures. Furthermore, we have evaluated our model on speech intelligibility with human listeners. Here, we use a sentence database that contains short semantically predictable (or natural) sentences for intelligibility tests. Figure 3 gives the keyword intelligibility score (median values and interquartile ranges) for the two-source configuration. Three SNR levels are tested: 0 dB, -5 dB and -10 dB, where the SNR is computed at the better ear, i.e. the one with higher SNR, for each sentence. The interfering source used for this configuration is babble noise. The general finding is that our algorithm improves human intelligibility for the tested conditions. The improvement becomes larger as the SNR decreases (61% at -10 dB), even though the algorithm introduces more target distortions at lower SNR levels. Figure 3b

shows the results for a three-source configuration, where our model yields a 40% improvement. Here the SNR is fixed at -10 dB at the better ear. The two interfering sources are one female speaker and a different male speaker. As far as we know, our system is the first binaural model that has been shown to produce a large speech intelligibility improvement for normal listeners.

A short paper describing this work was presented in 2002 *ICASSP*, and 2003 *NIPS*. A comprehensive version has been recently published by *Journal of the Acoustical Society of America*. Detailed references are given in Section 2.3 (Roman, Wang, and Brown, 2002b; 2003; 2004). Also, a powerpoint presentation file - titled "Speech segregation based on sound localization" - is posted on the PI's webpage at <http://www.cis.ohio-state.edu/~dwang> (follow the Presentations link), and the file contains sound demos. In addition, the source program that implements the algorithm is posted at the PI's laboratory webpage at <http://www.cis.ohio-state.edu/pnl> (follow the Software link).



(a)



(b)

**Figure 3.** Speech intelligibility score (median values and interquartile ranges) before (white bars) and after processing (black bars). (a) A two-source condition ( $0^\circ$  and  $5^\circ$ ) at three SNR levels: 0 dB, -5 dB and -10 dB. (b) A three-source condition ( $0^\circ$ ,  $30^\circ$  and  $-30^\circ$ ) at -10 dB SNR. In both conditions, target is at  $0^\circ$ .

## **2. OTHER INFORMATION**

### **2.1 Development of Human Resources**

Three doctoral students have been supported under this grant: Mingyang Wu, Nicoleta Roman, and Guoning Hu. Wu's research on multipitch tracking led to a Ph.D. dissertation successfully defended in September 2003. His dissertation, entitled "Pitch tracking and speech enhancement in noisy and reverberant environments", will soon be posted on the PI's laboratory webpage at <http://www.cis.ohio-state.edu/pnl>. An executive summary of the dissertation is given in Appendix 1.

Roman studies location-based speech segregation, and she has completed her Ph.D. Candidacy Examination, and is in the final stage of graduate study. Hu's work is on pitch-based speech segregation, and he has also completed his Candidacy Examination.

This grant has helped the PI to develop a graduate-level course entitled "Computational audition", and enhance the existing graduate-level courses "Survey of Artificial Intelligence", "Introduction to Artificial Neural Networks" and "Brain Theory and Neural Networks". Additionally, the PI has participated in a great deal of curriculum and seminar activity for training undergraduate students.

### **2.2 Transition or Collaborative Activities**

The PI visited AFRL in Dayton, Ohio, in January 2001. The visit was hosted by Dr. Tim Anderson, and the PI gave a presentation on the work performed in his group and explained the research objectives of the AFOSR project. In May 2001, Dr. Douglas Brungart of AFRL (Dayton OH) visited the PI's laboratory and discussed potential topics for collaboration. The continued interaction with Dr. Brungart finally led to a collaborative project on psychoacoustic evaluation of computational multitalker analysis systems since January 2003, and the project has revealed new insight into energetic and informational masking in human listeners by employing the concept of ideal binary masking originated in the PI's laboratory.

In October 2001, the PI interacted with Dr. Mark Ericson of AFRL (Dayton OH) at a conference, and in May 2002, Dr. Ericson visited the PI's laboratory and discussed potential topics for collaboration on auditory motion perception.

The PI is currently seeking collaboration with the AFRL/IF in Rome, New York. This AFRL facility is interested in speaker recognition in co-channel conditions (i.e. two speakers recorded on a single communication channel), and the PI plans to apply the results from this AFOSR project to solve the co-channel speaker recognition problem.

The PI worked closely with Dr. Willard Larkin of AFOSR in organizing the Workshop on Computational Audition, held during August 9-10, 2002, on OSU campus. The workshop featured thirteen speakers, who are leading experts in computational audition, and was attended by 40 participants from AFRL, OSU, and other institutions (e.g. the University of Illinois). In addition to presentations, the workshop generated a lot of in-depth discussion on related issues. Judging from the feedback of participants, the workshop was a considerable success.

The multipitch tracking algorithm described in Section 1.1 was included in a Matlab package for auditory research. This package is distributed free of charge to researchers



and practitioners around the world. A recent Ph.D. dissertation from the University of Plymouth in the U.K. made an independent comparison and confirmed the superior performance of this algorithm. The PI's laboratory also received many requests for software implementing the voiced speech segregation model and the binaural speech segregation model, described in Section 1. To facilitate technology transfer and benefit the broad research and development community, the laboratory posted software packages for multipitch tracking, voiced speech segregation, and binaural speech segregation on the worldwide web, which can be freely accessed on the internet.

Action Technologies, a small-business company located in Columbus, Ohio, partnered with the PI in winning a Phase I STTR project funded by AFOSR. This Phase I project seeks a feasibility study, to determine what improvements are needed in order to apply the monaural segregation algorithms developed in the PI's laboratory in practical situations. This project started in September 2003.

The PI provided consultation to a Canadian project on designing intelligent hearing aids. This joint project between the McMaster University and the Gennum Corporation, the world's largest hearing aids supplier, explores auditory scene analysis principles to improve current hearing aids technology.

## **2.3 Publications**

### **Journal articles:**

Van der Kouwe A.J.W., Wang D.L., and Brown G.J. (2001): "A comparison of auditory and blind separation techniques for speech segregation," *IEEE Transactions on Speech and Audio Processing*, vol. 9, pp. 189-195.

Wang D.L. and Liu X. (2002): "Scene analysis by integrating primitive segmentation and associative memory," *IEEE Transactions on Systems, Man, and Cybernetics, Part B: Cybernetics*, vol. 32, pp. 254-268.

Wu M., Wang D.L., and Brown G.J. (2003): "A multi-pitch tracking algorithm for noisy speech," *IEEE Transactions on Speech and Audio Processing*, vol. 11, 229-241.

Roman N., Wang D.L., Brown G.J. (2003): "Speech segregation based on sound localization," *Journal of the Acoustical Society of America*, vol. 114, pp. 2236-2252.

Hu G. and Wang D.L. (2004): "Monaural speech segregation based on pitch tracking and amplitude modulation," *IEEE Transactions on Neural Networks*, in press.

### **Book chapters:**

Wang D.L. (2003): "Temporal pattern processing," In: Arbib M.A. (ed.), *The Handbook of Brain Theory and Neural Networks*, 2nd Ed., MIT Press, Cambridge MA, pp. 1163-1167.

Brown G.J. and Wang D.L. (2004): "Timing is of the essence: Neural oscillator models of auditory grouping," In: Greenberg S. and Ainsworth W. (ed.), *Listening to Speech: An Auditory Perspective*, Lawrence Erlbaum, Mahwah NJ, in press.

#### Conference papers:

Wu M., Wang D.L., and Brown G.J. (2001): "Pitch tracking based on statistical anticipation," *Proceedings of International Joint Conference on Neural Networks (IJCNN-01)*, pp. 866-871.

Hu G. and Wang D.L. (2001): "An extended model for speech segregation," *Proceedings of IJCNN-01*, pp. 1089-1094.

Roman N., Wang D.L., and Brown G.J. (2001): "Speech segregation based on sound localization," *Proceedings of IJCNN-01*, pp. 2861-2866.

Brown G.J., Barker J., and Wang D.L. (2001): "A neural oscillator sound separator for missing data speech recognition," *Proceedings of IJCNN-01*, pp. 2907-2912.

Palomaki K., Brown G.J., and Wang D.L. (2001): "A binaural model for missing data speech recognition in noisy and reverberant conditions," *Web Proceedings of Workshop on Consistent and Reliable Acoustic Cues for Sound Analysis*.

Hu G. and Wang D.L. (2001): "Speech segregation based on pitch tracking and amplitude modulation," *Proceedings of IEEE Workshop on Applications of Signal Processing to Audio and Acoustics*, pp. 79-82.

Wu M., Wang D.L., and Brown G.J. (2002): "A multi-pitch tracking algorithm for noisy speech," *Proceedings of IEEE International Conference on Acoustics, Speech, & Signal Processing (ICASSP)*, pp. 369-372.

Hu G. and Wang D.L. (2002a): "On amplitude modulation for monaural speech segregation," *Proceedings of IJCNN-02*, pp. 69-74.

Hu G. and Wang D.L. (2002b): "Monaural speech segregation based on pitch tracking and amplitude modulation," *Proceedings of ICASSP-02*, pp. 553-556.

Roman N., Wang D.L., and Brown G.J. (2002a): "Localization-based sound segregation," *Proceedings of IJCNN-02*, pp. 2299-2303.

Roman N., Wang D.L., and Brown G.J. (2002b): "Localization-based sound segregation," *Proceedings of ICASSP-02*, pp. 1013-1016.

Hu G. and Wang D.L. (2003): "Monaural speech separation," in *Advances in Neural Information Processing Systems (NIPS-02)*, vol. 15, pp. 1221-1228, Cambridge MA: MIT Press, 2003.

Hu G. and Wang D.L. (2003): "Separation of stop consonants," *Proceedings of ICASSP-03*, pp. II.749-752.

Roman N. and Wang D.L. (2003): "Binaural tracking of multiple moving sources," *Proceedings of ICASSP-03*, pp. V.149-152.

Shao Y. and Wang D.L. (2003): "Co-channel speaker identification using usable speech extraction based on multi-pitch tracking," *Proceedings of ICASSP-03*, vol. II.205-208.

Wu M. and Wang D.L. (2003): "A one-microphone algorithm for reverberant speech enhancement," *Proceedings of ICASSP-03*, pp. I.844-847.

Srinivasan S. and Wang D.L. (2003): "Schema-based modeling of phonemic restoration," *Proceedings of EUROSPEECH-03*, pp. 2053-2056.

Roman N., Wang D.L., and Brown G.J. (2004): "A classification-based cocktail-party processor," *Proceedings of NIPS-03*, in press.

Roman N. and Wang D.L. (2004): "Binaural sound segregation for multisource reverberant environments," in *Proceedings of ICASSP-04*, to appear.

# **Report of Inventions and Subcontracts**

DeLiang Wang  
(Principal Investigator)

March 2004

*Department of Computer and Information Science and Center for Cognitive Science  
The Ohio State University*

The project, entitled "Computational Auditory Scene Analysis Based Perceptual and Neural Principles" (F49620-01-1-0027), was funded by the Air Force Office of Scientific Research from December 2000 through November 2003.

This is to certify that no patent was filed or planned, as a result of the above AFOSR grant.

## **Appendix 1. Executive Summary of Mingyang Wu's Ph.D. Dissertation**

Two causes of speech degradation exist in practically all listening situations: noise interference and room reverberation. This dissertation investigates three particular aspects of speech processing in noisy and reverberant environments: multipitch tracking for noisy speech, measurement of reverberation time based on pitch strength, and reverberant speech enhancement using one microphone (or monaurally).

An effective multipitch tracking algorithm for noisy speech is critical for speech analysis and processing. We present a robust algorithm for multipitch tracking of noisy speech. Our approach integrates an improved channel and peak selection method, a new method for extracting periodicity information across different channels, and a hidden Markov model (HMM) for forming continuous pitch tracks. The resulting algorithm can reliably track single and double pitch tracks in a noisy environment. We evaluate our algorithm on a database of speech utterances mixed with various types of interference. Quantitative comparisons show that our algorithm significantly outperforms existing ones.

Reverberation corrupts harmonic structure in voiced speech. We observe that the pitch strength of voiced speech segments is indicative of the degree of reverberation. Consequently, we present a pitch-based measure for reverberation time (T60) utilizing our new pitch determination algorithm. The pitch strength is measured by deriving the statistics of relative time lags, defined as the distances from the detected pitch periods to the closest peaks in correlograms. The monotonic relationship between the measured pitch strength and reverberation time is learned from a corpus of reverberant speech with known reverberation times.

Under noise-free conditions, the quality of reverberant speech is dependent on two distinct perceptual components: coloration and long-term reverberation. They correspond to two physical variables: signal-to-reverberant energy ratio (SRR) and reverberation time, respectively. We propose a two-stage reverberant speech enhancement algorithm using one microphone. In the first stage, an inverse filter is estimated to reduce coloration effects so that SRR is increased. The second stage utilizes spectral subtraction to minimize the influence of long-term reverberation. The proposed algorithm significantly improves the quality of reverberant speech. Our algorithm is quantitatively compared with a recent one-microphone reverberant speech enhancement algorithm. The results show that our algorithm performs substantially better.